# *DUOS:*
# *A Structured Approach to Genomics' Data Use Oversight*

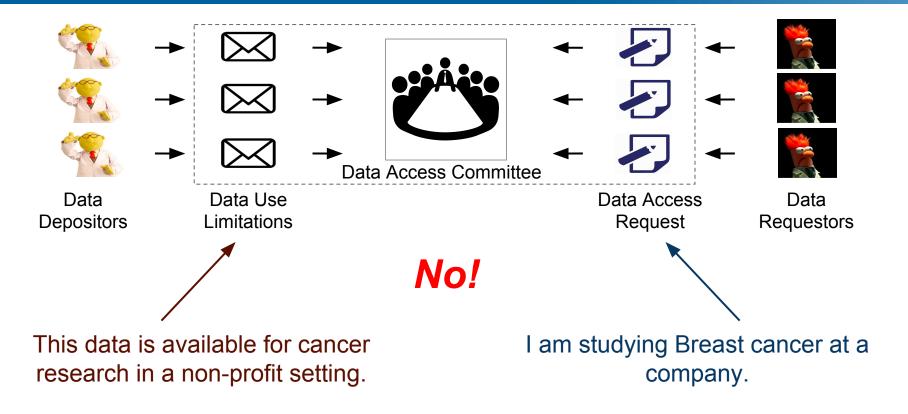https://duos.broadinstitute.org

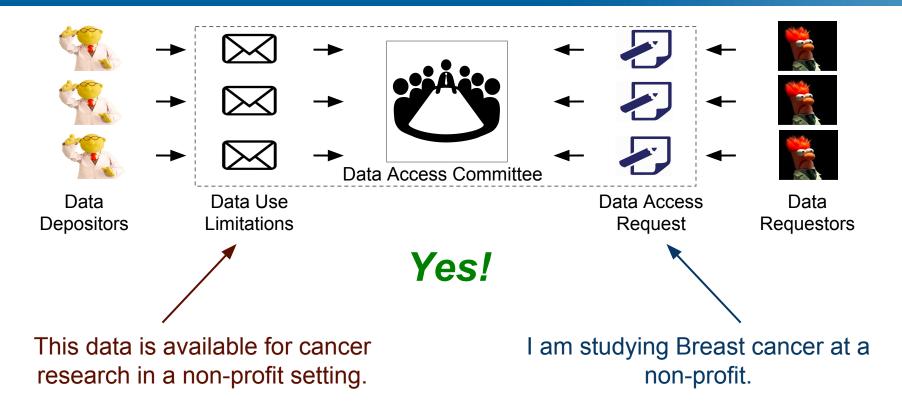**BROAD** INSTITUTE

*Moran Cabili*
*nmcabili@broadinstitute.org*

# If Time From Discovery to Access Was Instant...

# Our Current Protocol for Data Access

Data Access Committee

Data Depositors

Data Use Limitations

Data Access Request

Data Requestors

*No!*

This data is available for cancer research in a non-profit setting.

I am studying Breast cancer at a company.

Data Access Committee

Data
Depositors

Data Use
Limitations

Data Access
Request

Data
Requestors

*Yes!*

This data is available for cancer
research in a non-profit setting.

I am studying Breast cancer at a
non-profit.

# Our Current Protocol for Data Access



Data Access Committee

Data Depositors → Data Use Limitations → Data Access Request ← Data Requestors

**Scales Poorly!!**
**$O(N^2)$**

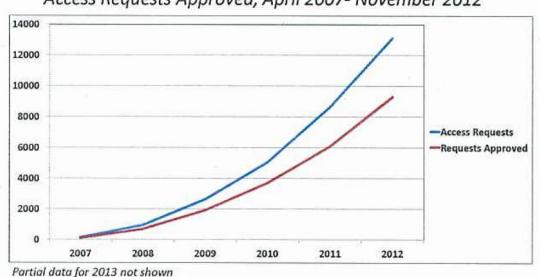dbGaP at PRIMr 2013

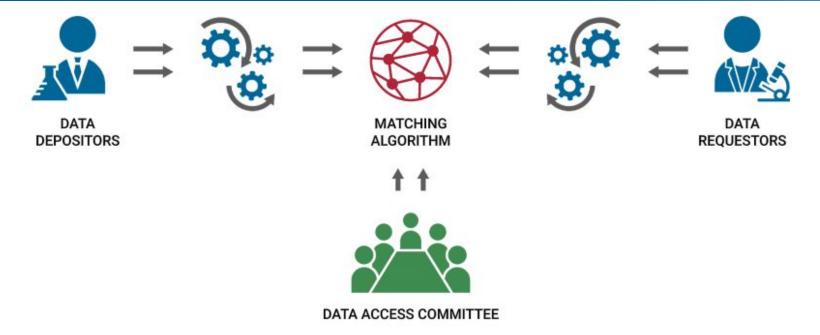Cumulative Number of Data Access Requests Submitted and Data Access Requests Approved, April 2007- November 2012

Access Requests
Requests Approved

Partial data for 2013 not shown

# DUOS



## *What is DUOS?*

- Interfaces to transform data use restrictions and data access requests to machine-readable code (ADA-M & Consent Codes)

- A matching algorithm that checks if data access requests are compatible with data use restrictions

- Interfaces for the Data Access Committee to adjudicate whether structuring and matching has been done appropriately
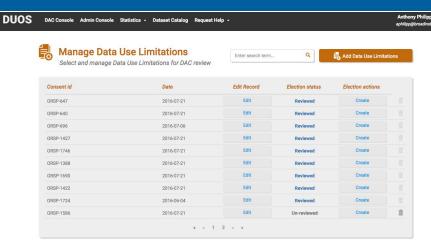
# DUOS: Matching Algorithm
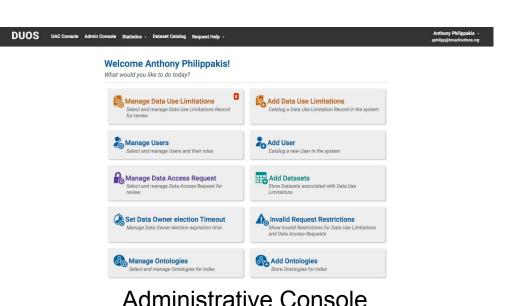
## Uses a Web Ontology Language (OWL) reasoner

### Consent Category

| | |
|---|---|
| Restricted to the study of Disease (ontology ID) | DS-[XX](CC) **or** HMB(CC) |
| Commercial use is prohibited | NPU |
| General methods development prohibited | NMDS |
| Population structure or normal variation studies prohibited | NPNV |
| Samples use as controls outside the specified restrictions is prohibited | *NCTRL* |
| limited to gender (M/F/none) | RS-[GENDER] |
| limited to pediatric research | RS-[PEDIATRIC] |
| Sensitive/restricted population; Date | MAN-REV |

Data use restrictions  VS  Data access request

Disease — Cancer — Breast Cancer

Disease — Cancer — Breast Cancer

DUL defaults

Study males

For commercial use

## Matching Algorithm
## Answer:  *Yes!*

# DUOS: DAC Interfaces



Structure & Manage Data Use Letters

Approval of Structured Data Use

Administrative Console

DAC Members Approve Access

# Validation of DUOS

## **_Claim:_** Data Use Can Be Structured



Figure 2: Data use restrictions can be structured by using 5 main categories (non-mutually exclusive; N=125)

**Diseases:** Diabetes research only, Breast cancer research only, etc
**Commercial Use:** allowed/not allowed.
**Special populations:** Ethnicities, gender, pediatric, etc.
**Future use for Methods Development, Aggregate Statistics, Controls**

*Review of ~150 Data Use Limitations Letters at Broad demonstrated that ~90% can be structured with the following ontologies*

## **_Test:_** Run a trial!

**Data Access Committee**
Pearl O'Rourke (Partners)
Laura Rodriguez (NIH)
John Wilbanks (Sage)
Stacey Donnelly (Broad)
Anthony Philippakis (Broad)

*We have formed a DAC to compare automated review of access to traditional mode.*

# Trial of Structured Data Use



Data Depositors | Data Use Limitations | Data Access Committee | Project Request Forms | Data Requestors

1     3     2

1) Can data use letters be structured?

2) Can research purpose be structured?

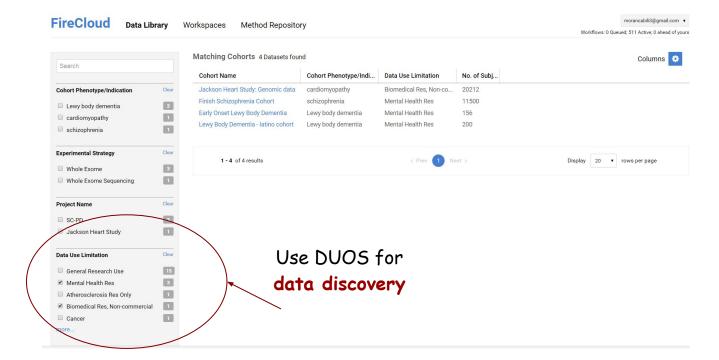3) Does the DAC agree with the verdict of DUOS?

# What's next?

1. Updating our ontology to reference DUO
   a. Propose adding new categories to DUO (adopted from ADA-M)

2. DUOS pilot: 3rd batch evaluation of structuring Data Use Limitations

3. Connect to a cloud data repository:  Broad's Data Library on 



Use DUOS for **data discovery**

# Acknowledgements

Greg Rushton
Andrea Saltzman
Stacey Donnelly
Anthony Philippakis
Moran Cabili

Pearl O'Rourke
Laura Rodriguez
John Wilbanks

Broad Data Science Platform
DSDE dev team

Stephanie OM Dyke
Anthony Brooks
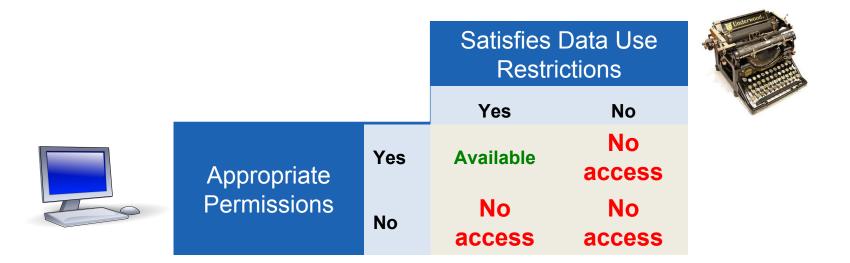Melanie Courtot
Dylan Spalding

# DUOS Ontology is Consistent with DUO and ADA-M

| DUOS Restriction code | Dyke et al. Consent Code | abbreviation | ADA-M |
|---|---|---|---|
| **Primary Codes: select one** | | | |
| No restrictions | no restrictions | NRES | v |
| General research use | general research use and clinical care | GRU(CC) | v |
| Health and Biomedical Research | health/medical/biomedical research and clinical care | HMB(CC) | v |
| Restricted to the study of Disease (ontology ID) | disease-specific research and clinical care | DS-[XX](CC) | v |
| population origins/ancestry research | population origins/ancestry research | POA | v |
| **Secondary codes: select any** | | | |
| Commercial use is prohibited | not-for-profit use only | NPU | v |
| General methods development prohibited | no "general methods" research | NMDS | v |
| Population structure or normal variation studies prohibited (NPNV) *(\*dbGaP question on submission forms)* | N/A | NPNV | v(2 categories) |
| Samples use as controls outside the specified restrictions is prohibited(NCTRL) | NOT MODELED EXPLICITLY; will be set to True if the following primary categories were selected: disease-specific research and clinical care OR population origins/ancestry research | DS-[XX](CC) OR POA | v |
| Future use of aggregate-level data for general research purposes is prohibited [NAGR] *(\*dbGaP question on submission forms)* | N/A | NAGR | N/A |
| limited to gender (M/F/none) | other research-specific restrictions | RS-[GENDER] | v |
| limited to pediatric research | other research-specific restrictions | RS-[PEDIATRIC] | v |
| Sensitive/restricted population; Date | | | |
| N/A | research use only | RUO | v |
| N/A | genetic studies only | GSO | v |
| Recontacting Data Subjects (*may* and *must*) | | | v |
| Other term of use (ethics committee req, cloud storage prohibited, geographical restriction) | ethics approval required,geographical restrictions | IRB, GS[XX] | v |

# Problem: Data Use is not Coded!

***Data Use Restrictions****: What are you doing with the data?*

"The donor wants her data used only for non-commercial cancer research"

|  | | Satisfies Data Use Restrictions | |
|---|---|---|---|
|  | | Yes | No |
| Appropriate Permissions | Yes | Available | No access |
|  | No | No access | No access |

***Permissions****: Who are you?*

"Only consortium members can READ this data until it is published."

## Main Question: Can Data Use Restrictions be made machine-readable?